# REANA and ML

Tibor Simko
IT-PW

CERN IT Machine Learning Infrastructure Workshop, March 10th 2023
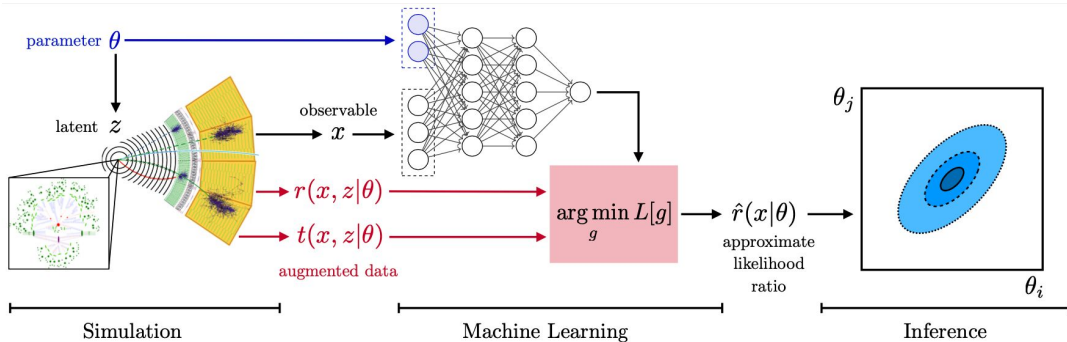
https://indico.cern.ch/event/1253881

# REANA Reusable Analysis platform

Running declarative containerised computational workflows

# ML use cases on REANA 1/2

Pheno-level analyses embedded into Python ML ecosystem (and optionally MLFlow)
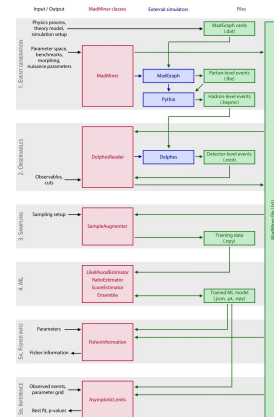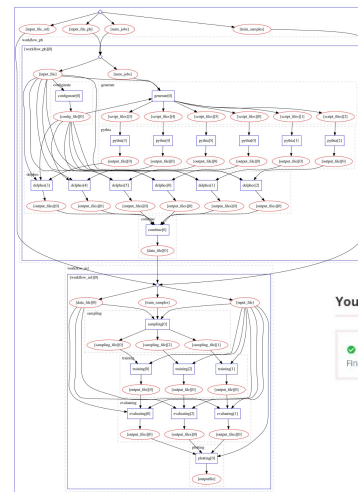




"MadMiner: Machine learning-based inference for particle physics", J. Brehmer, F. Kling, I. Espejo, K. Cranmer, [arXiv:1907.10621](arXiv:1907.10621).

- Running ML based workflows

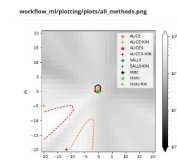# ML use cases on REANA 2/2



"REANA / PanDA integration for Active Learning", W.Guan, T.Maeno,
C.Weber, T.Wenaus, R.Zhang, https://indico.cern.ch/event/1134581.

ATL-PHYS-PUB-2022-045

- Running workflows as part of a bigger data processing chain (whole physics
  analysis from MC generation to new physics discovery)

# Possible areas of interest

- Capturing the knowledge behind data analyses
  - → *preserve to reuse*
- Computational reproducibility
  - → *run outside the original context*
- Running workflows at scale
  - → *10k workflows for ATLAS pMSSM searches*
- "Continuous analyses"
  - → *Gitlab-REANA bridge*
- Interplay between notebooks and workflows
  - → *interactive vs batch*

A Large-scale Study about Quality and Reproducibility of Jupyter Notebooks

João Felipe Pimentel*, Leonardo Murta*, Vanessa Braganholo*, and Juliana Freire†
*Universidade Federal Fluminense
Niterói, Brazil
{jpimentel,leomurta,vanessa}@ic.uff.br
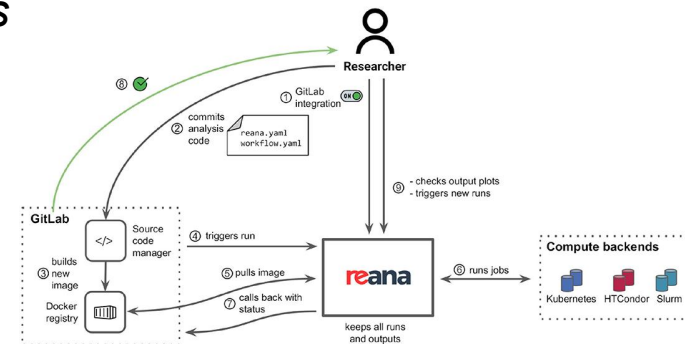†New York University
New York, USA
juliana.freire@nyu.edu

"...only 4.03% produced the same results"
DOI 10.1109/MSR.2019.00077



GitLab-REANA bridge